

خوارزمية عامة لتوسيع بيانات تخضع لتوزيع بيتا متعدد المتغيرات

د. صلاح حمزه عبد*

المستخلص

سيق في عام 2002 وقدم كل من Michael و Schucany [2] ، جهداً رائعاً ، تمثل في أسلوبهما الذي أطلقوا عليه تسمية الأسلوب الخلط لتوسيع بيانات تخضع لتوزيعات ثنائية المتغيرات Bivariate . سنقوم في هذا البحث بطرح خوارزمية عامة لتوسيع بيانات تخضع لتوزيع بيتا متعدد المتغيرات باستخدام التعميم الذي طرحته كل من Minhajuddin و Harris و Schucany عام 2004 [3] للأسلوب الخلط المنوه عنه أعلاه ، إذ جعلوه ممكناً التطبيق على التوزيعات متعددة المتغيرات .

١- المقدمة

لا يخفى على المتبع للنظرية الإحصائية أهمية توليد المشاهدات على وفق أنموذج أو مسالة ما للأغراض التجريبية . وعلى الرغم من أن هذا المسئلہ قد يكون صعباً في أحيان كثيرة ، إلا أنه يزداد صعوبة في حالة المتغيرات المتعددة ، بسبب الارتباط الحاصل ما بين المتغيرات التي تخضع بمجملها للتوزيع المتعدد المتغيرات ، وكيفية استكشافه ، واحتسابه والتعامل معه . إن الدراسات القائمة على المحاكاة تستند على توليد المشاهدات على وفق هذا التوزيع الاحتمالي أو ذلك لتحقيق جوانب رياضية نظرية قد تعجز الطرق التحليلية الرياضية عن تحقيقها بعض المسائل المتعلقة بأحد شقي النظرية الإحصائية ، التقدير أو الاختبار . إن استئناد السواد الأعظم من النظرية الإحصائية على فرض التوزيع الطبيعي ، سواءً أكان أحادي أو متعدد للمتغيرات ، لا يعني أن ليس هناك ظواهر في الواقع العملي لا تخضع لهذا التوزيع ، بل قد ثبتت الدراسات بأن معظم البيانات المأخوذة من ظواهر حقيقة ، هي تخضع لتوزيعات

* اسنان/قسم الإحصاء/كلية الإدارة والاقتصاد/جامعة المستنصرية
(15)

غير طبيعية [2] ، إلا أن محاكاة الظواهر ذات المتغيرات المتعددة التي لا تخضع للتوزيع الطبيعي ليست شائعة بسبب عدم توفر الخوارزميات التي تجعل من ذلك ممكناً [3] .

إن أحد ابرز التوزيعات شائعة الاستخدام في الجوانب النظرية والتطبيقية سواء أكان أحادي أو متعدد للمتغيرات ، يتمثل بتوزيع بيتا ، إذ يقال بأن للمتغير y توزيع بيتا بالمعلمتين a و b ،
إذا امتلك دالة كثافة الاحتمال :

$$f_y(y) = \frac{y^{a-1}(1-y)^{b-1}}{B(a,b)} , \quad 0 < y < 1 , \quad a, b > 0 , \quad \dots \dots \dots (1)$$

حيث أن $B(a,b) = \int_0^1 y^{a-1}(1-y)^{b-1} dy$ ، عبارة عن دالة بيتا الرياضية المعروفة .

أما لو افترضنا أن عينة عشوائية بحجم n قد سحب من مجتمع يخضع للتوزيع الطبيعي ذي Σ من المتغيرات ، بحيث أن $m \leq n$ ، بمتجه المتوسطات \bar{x} ومصفوفة تباين وتبالين مشترك Σ ، فإن مقرر مصفوفة التباين وتبالين المشترك :

$$S = \frac{\sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^T}{n-1} = \frac{A}{n-1} \quad \dots \dots \dots (2)$$

سيخضع للتوزيع وشرت (Wishart distribution) ، بـ $r = n - 1$ من درجات الحرية و
كمصفوفة متوسطات ، وعلى وفق دالة كثافة الاحتمال [1] ،

$$f_S(S) = \begin{cases} \frac{|S|^{\frac{r-m-1}{2}} e^{\frac{-1}{2} \text{tr } S^{-1} S}}{2^{\frac{rn}{2}} \pi^{\frac{m(m-1)}{4}} |V|^{\frac{r}{2}} \prod_{i=1}^m \left(\frac{r+1-i}{2} \right)} & , \text{ if } S \text{ is positive definite} \\ 0 & , \text{ o.w.} \end{cases} \quad \dots \dots \dots (3)$$

(16)

إذ يشار في الأدبيات العلمية لما ورد أعلاه بالشكل $s \sim W_m(r, V)$

فلو كان لدينا $(j=1,2)$ ، $s_j \sim W_m(r_j, V)$ ، حيث أن s_1 و s_2 مستقلتين عن بعضهما البعض ، فإنه سيقال بأن المصفوفة $H = C^{-1}S_2(C^{-1})^T$ توزع بيتا متعدد المتغيرات [1] ، وبذلة كثافة الاحتمال (Multivariate Beta distribution) :

$$f_H(h) = \begin{cases} \frac{\Gamma_m\left(\frac{r_1+r_2}{2}\right)}{\Gamma_m\left(\frac{r_1}{2}\right)\Gamma_m\left(\frac{r_2}{2}\right)} |h|^{(r_2-m-1)/2} |1-h|^{(r_1-m-1)/2}, & h>0, 1-h>0 \\ 0 & \text{o.w.} \end{cases} \quad (4)$$

حيث أن $C = S_1 + S_2$ ، وأن $C C' = S_1 + S_2$ هي مصفوفة مثلث سفلي .

"وبناء" على ما ورد أعلاه ، ولعدم توفر الخوارزميات [3] ، التي تجعل من محاكاة الظواهر المستندة على متغيرات متعددة أمراً ممكناً ، فقد آثرنا هنا تقديم خوارزمية عامة لتوليد بيانات تخصيص توزيع بيتا متعدد المتغيرات ، لممثل هدف بحثنا هذا .

The Mixture Method [2,3]

إذا افترضنا بأن للمتغير العشوائي K دالة كثافة احتمال (إذا كان المتغير K مستمراً) ، أو دالة كتلة احتمال (إذا كان لمتغير العشوائي K متقطعاً) ، تتمثل بالشكل :

$$u(k; \theta) \quad (5)$$

على أن θ يمكن أن تكون معلمة واحدة أو متوجه من المعالم ، فإنه بالاعتماد على $K = k$ يمكن توليد m من المتغيرات x_1, x_2, \dots, x_m ، المستقلة عن بعضها البعض ، وكل منها يمتلك دالة كثافة احتمال شرطية (أو دالة كتلة احتمال شرطية) ،

(17)

$$p(x_i | k, \eta) \quad , \quad i = 1, 2, \dots, m \quad \dots \dots \dots \quad (6)$$

حيث أن η يمكن أن تكون هي الأخرى معلمة واحدة أو متوجه من المعامل .

وبضرب المعادلين في (5) و (6) مع بعضهما البعض ، سنحصل على دالة كثافة (كتلة)

احتمال مشتركة لـ x_1, x_2, \dots, x_m و K ، بالشكل :

$$g(x_1, x_2, \dots, x_m, k, \theta, \eta) = u(k, \theta) \prod_{i=1}^m p(x_i | k, \eta) \quad \dots \dots \dots \quad (7)$$

نقد لاحظ **Minhajuddin** و **Harris** و **Schucany** عام 2004 [3] ، بأن دالة كثافة (أو كتلة) الاحتمال المشتركة الواردة في المعادلة (7) متماثلة (symmetric) بالنسبة لأي متغيرين مثل x_i و x_j ، حيث أن $i \neq j$ ، وأن $i, j = 1, 2, \dots, m$ ، فاستنتجوا بأن المتغيرات العشوائية x_1, x_2, \dots, x_m ، قابلة للتبادل (Exchangable) ، وأنه بالاعتماد على $K = k$ ، فإن x_i و x_j سيكونا مستقلين ، وأنه يمكن الحصول على دالة كثافة (أو كتلة) الاحتمال المشتركة للمتغيرات x_1, x_2, \dots, x_m بالشكل ،

$$f(x_1, x_2, \dots, x_m; \theta, \eta) = \begin{cases} \int g(x_1, x_2, \dots, x_m, k; \theta, \eta) dk , & \text{if } K \text{ continuous} \\ \sum_k g(x_1, x_2, \dots, x_m, k; \theta, \eta) , & \text{if } K \text{ discrete} \end{cases} \quad \dots \dots \dots \quad (8)$$

ان قيمة المعلمة (أو قيمة متوجه المعامل) η في الدالة $f(x_1, x_2, \dots, x_m; \theta, \eta)$ ، تلعب دوراً "مهماً" جداً في السيطرة على الارتباط ما بين المتغيرات القابلة للتبادل x_i و x_j ($i \neq j$) ،

على وفق أقيام المشاهدات التي تخضع للمتغير K .

وببناءً على الفلسفة الواردة في أعلى ، خلص كل من **Minhajuddin** و **Harris** و **Schucany** عام 2004 [3] ، إلى الخطوتين التاليتين لمحاكاة التوزيعات ذات المتغيرات المتعددة.

(18)

(1) محاكاة المتغير $K = k$ على وفق دالة كثافة (كتلة) الاحتمال الحدية $u(k; \theta)$

(2) اعتماداً على كل قيمة مولدة للمتغير $k = K$ ، يتم توليد m من القيم المستقلة للمتغيرات x_i ($i=1,2,\dots,m$) ، قيمة لكل متغير ، وذلك من خلال التوزيع اللاحق $p(x_i | K = k, \eta)$

إن أحد أبرز إفرازات نظرية بيز ، قد تمثل بعوامل التوزيعات الأولية المترافق ، فعائلة توزيع بيتا مثلاً هي المترافق لعائلة توزيع ذي الحدين ، بمعنى أنه لو كان X يخضع لتوزيع بيتا بالعلمين α و β ، وأن $X = K$ يخضع لتوزيع ذي الحدين ، فإن هذا سيؤدي إلى خصوص K لتوزيع حدي (غير معتمد) هو عبارة عن توزيع بيتا - ذي الحدين (Beta-Binomial) ، كما أن عائلة توزيع كاما هي المترافق لعائلة توزيع بواسون ، وهذا يؤدي إلى خصوص K لتوزيع حدي (غير معتمد) هو عبارة عن توزيع ذي الحدين السالب .

لقد استفاد منظروا الطريقة الخليطة قيد الدراسة من إفراز نظرية بيز المذكور أعلاه ، من خلال البدء بالتوزيع الحدي (غير المعتمد) للمتغير K ، وبالاعتماد على القيمة المولدة $K = k$ ، تم توليد أقيام المتغيرات x_1, x_2, \dots, x_m ، وكل على وفق العائلة المترافق معها .

III- محاكاة عائلة بيتا متعددة المتغيرات

يذكر كل من Patil و Boswell و Joshi و Ratnaparkhi عام 1984 [5] ، الصيغة

التالية لتوزيع بيتا - ذي الحدين :

$$p(k) = \frac{C_k^{-\alpha} C_{v-k}^{-\beta}}{C_v^{-\alpha-\beta}} , \quad k = 0, 1, \dots, v , \quad v = 1, 2, \dots \quad --- (9)$$

ويستخدم العلاقة الرياضية [4] ، فإنه يمكن كتابة $p(k)$

في المعادلة (9) بالشكل :

(19)

$$p(k) = \frac{(-1)^k C_k^{\alpha+k-1} (-1)^{v-k} C_{v-k}^{v-k+\beta-1}}{(-1)^v C_v^{v+\alpha+\beta-1}}$$

وبالاختصار وإعادة كتابة التوافق بشكل مفهوم ، يكون ،

$$\begin{aligned} p(k) &= \frac{v!(\alpha+k-1)!(v-k+\beta-1)!(\alpha+\beta-1)!}{k!(\alpha-1)!(\beta-1)!(v-k)!(v+\alpha+\beta-1)!} \\ &= C_k^v \frac{\Gamma(k+\alpha)\Gamma(v+\beta-k)\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)\Gamma(v+\alpha+\beta)} \\ &= C_k^v \frac{B(k+v, v+\beta-k)}{B(\alpha, \beta)} \quad --- (10) \end{aligned}$$

والصيغة المذكورة في (10) أعلاه هي شكل آخر لدالة كثافة احتمال توزيع بيتا - ذي الحدين .

إن متوسط وتباعد المتغير K سيكونا عباره عن [5]

$$E(K) = \frac{v\alpha}{\alpha + \beta}, \quad Var(K) = \frac{v\alpha\beta(\alpha + \beta + v)}{(\alpha + \beta)^2(\alpha + \beta + 1)} \quad --- (11)$$

ومن ذلك فأن ،

$$E(K^2) = Var(K) + [E(K)]^2 = \frac{v\alpha(v\alpha + v + \beta)}{(\alpha + \beta + 1)(\alpha + \beta)} \quad --- (12)$$

وتلبيساً على ما ورد أعلاه وما ورد في المعادلة (7) ، فان دالة الكثافة الاحتمالية المشتركة مابين أي من المتغيرات x_i ($i=1, 2, \dots, m$) والمتغير K ستكون :

$$\begin{aligned} g(x_i, k; v, \alpha, \beta) &= \frac{x_i^{\alpha+k-1} (1-x_i)^{v-\beta-k-1}}{B(k+\alpha, v+\beta-k)} \cdot \frac{C_k^v B(k+\alpha, v+\beta-k)}{B(\alpha, \beta)} \\ &= \frac{x_i^{\alpha-1} (1-x_i)^{\beta-1}}{B(\alpha, \beta)} \cdot C_k^v x_i^k (1-x_i)^{v-k} \quad --- (13) \end{aligned}$$

إذا أن $x_i | k$ ($i=1, 2, \dots, m$) ، يخضع توزيع بيتا بالمعلمتين $v + \beta - k$ و $\alpha + k$ (20)

وبأخذ المجموع لكل قيم K ، حيث أن $v = 0, 1, 2, \dots, n$ ، يمكن الحصول على دالة الاحتمال الحدية للمتغير x_i ، إذ أن الحدود المتضمنة المتغير K ، في المعادلة (13) ، عبارة عن دالة ذي الحدين بالمعلمتين v و x_i ، ومجموعها يؤول للواحد الصحيح ، فيكون توزيع أي متغير مثل $(x_i)_{i=1,2,\dots,m}$ ، عبارة عن توزيع بيتا بالمعلمتين α و β ، وعليه فإن متوسط وتبان أي متغير مثل x_i سيكون عبارة عن :

$$E(x_i) = \frac{\alpha}{\alpha + \beta} , \quad Var(x_i) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} \quad \text{---(14)}$$

وبما أنه يتم توليد قيمة لكل متغير من المتغيرات x_1, x_2, \dots, x_m بشكل مستقل عن القيمة المولدة لمتغير آخر ، اعتماداً على نفس القيمة للمتغير K ، كما أسلفنا ، فـ "أنتادا" لمعادلتين (7) و (13) ، يكون التوزيع الشرطي لكل متغير مثل x_i ($i = 1, 2, \dots, m$) اعتماداً على $K = k$ ، هو عبارة عن توزيع بيتا بالمعلمتين $v + \beta - k$ و $\alpha + k$ ، وبدالة كثافة الاحتمال :

$$p(x_i | k) = \frac{x_i^{k+\alpha-1} (1-x_i)^{v+\beta-k-1}}{B(k+\alpha, v+\beta-k)} , \quad i = 1, 2, \dots, m \quad \text{---(15)}$$

وعلى ذلك فان :

$$\left. \begin{aligned} E_{X_i|K}(X_i) &= \frac{K+\alpha}{\alpha+\beta+v} \\ Var_{X_i|K}(X_i) &= \frac{(K+\alpha)(\beta+v-K)}{(\alpha+\beta+v+1)(\alpha+\beta+v)^2} \end{aligned} \right\} \quad \text{---(16)}$$

وأنه لكل $1 \leq i \neq j \leq m$ ، يكون :

$$\begin{aligned} E(X_i X_j) &= E_K [E_{X_i X_j | K}(X_i X_j)] \\ &= E_K [E_{X_i | K}(X_i) \cdot E_{X_j | K}(X_j)] \\ &= E_K \left[\frac{(K+\alpha)^2}{(\alpha+\beta+v)^2} \right] \\ &= E_K \left[\frac{K^2 + 2\alpha K + \alpha^2}{(\alpha+\beta+v)^2} \right] \end{aligned} \quad (21)$$

وبالتعويض من المعادلين (11) و (12) ، وجملة من العمليات البسيطة ، يكون :

$$E(X_i X_j) = \frac{\alpha \{v(\alpha+1) + \alpha(\alpha+\beta+1)\}}{(\alpha+\beta+v)(\alpha+\beta+1)(\alpha+\beta)}$$

ف تكون قيمة الارتباط ما بين أي متغيرين مثل x_i و x_j ($i \neq j$) عبارة عن :

$$\begin{aligned} \rho_{X_i X_j} &= \frac{E(X_i X_j) - E(X_i) E(X_j)}{\sqrt{Var(X_i) Var(X_j)}} \\ &= \frac{v\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)(\alpha+\beta+v)} \cdot \frac{(\alpha+\beta)^2(\alpha+\beta+1)}{\alpha\beta} \\ &= \frac{v}{(\alpha+\beta+v)} \end{aligned} \quad -(17)$$

إن قيمة الارتباط في (17) أعلاه ، ذات أهمية بالغة في عمل الخوارزمية العامة قيد الطرح ، إذ من خلالها يمكن تحديد قيمة المعلمة v في توزيع بيتا - ذي الحدين ، حيث يكون :

$$v = \frac{(\alpha+\beta)\rho}{1-\rho} \quad -(18)$$

فتدخل قيمة الارتباط ما بين المتغيرات في عمل الخوارزمية على وفق هذا النحو ، لما للارتباط من أهمية بالغة في تحليل العلاقات ما بين المتغيرات المتعددة .

IV. خوارزمية المحاكاة وتتفيد منها

بعد ما ورد في فقرات البحث السابقة ، وعلى الأخص الفقرة (III) ، يمكن تلخيص خوارزمية المحاكاة ، إجمالاً بالآتي :

(a) يتم توليد قيمة p ، لتخضع لتوزيع بيتا بالمعلمتين α و β ، وعلى وفق دالة كثافة الاحتمال في (1) ، إذ أن :

$$\beta = b \quad \text{و} \quad \alpha = a$$

(22)

(b) توليد قيمة k من خلال خوارزمية توزيع ذي الحدين ، على اعتبار قيمة p المولدة في الخطوة (a) أعلاه كقيمة معلمة لهذا التوزيع جنباً إلى جنب مع قيمة v المستخرجة من المعادلة (18) ، كقيمة معلمة هي الأخرى ، لتخضع k بالنتيجة لتوزيع بيتاً - ذي الحدين بالمعلمات v و $\alpha^* = \beta^*$ و $\beta^* = \alpha + \beta$

(c) توليد m من القيم x_1, x_2, \dots, x_m ، على أساس قيمة k المولدة في الخطوة (b) السابقة ، وذلك على وفق توزيع بيتاً بالمعلمتين $v + \beta^* - k$ و $\alpha^* + k$ ، لتدرج كل قيمة تحت المتغير المناظر لنرميزها X_1, X_2, \dots, X_m على التوالي كقيمة مولدة على وفقه .

(d) يتم تكرار الخطوات السابقة بقدر عدد المشاهدات المراد توليده .
أما البرنامج المكتوب من قبل الباحث بلغة فجوال بيذك الذي يعبر عن الخوارزمية أعلاه و يجعلها قابلة للتنفيذ فهو الآتي :

```

1 REM " THIS PROGRAM WRITTEN BY (DR. SALAH HAMZA
ABID) "
2 REM " TO GENERATE OBSERVATIONS FELLOW
MULTIVARIATE "
3 REM " BETA DISTRIBUTION WITH CORRELATION
COEFFICIENT "
4 REM " (rw) AND (m) OF VARIABLES AND (n) OF
OBSERVATIONS "
5 REM "
*****
10 INPUT "rw=";RW
20 INPUT "m=";M
30 INPUT "n=";N
40 INPUT "a=";A
50 INPUT "b=";B

```

```

60 DIM X(M,N)
70 FOR J=1 TO N:RANDOMIZE TIMER:SSA=0:SSB=0:k=0
80 FOR I=1 TO INT(A)
90 SSA=SSA-LOG(RND):NEXT
100 SSA=SSA-(LOG(RND))*(A-INT(A))
110 FOR I=1 TO INT(B):SSB=SSB-LOG(RND):NEXT
120 SSB=SSB-(LOG(RND))*(B-INT(B))
130 P=SSA/(SSA+SSB)
140 V=((A+B)*RW)/(1-RW)
150 FOR I=1 TO INT(V)
160 IF RND <= P THEN K=K+1
170 NEXT
180 IF ((V-INT(V))*RND) <= P THEN K=K+1
190 FOR I=1 TO M:SS1=0:SS2=0
200 FOR L1=1 TO INT(K+A+B):SS1=SS1-LOG(RND):NEXT
210 SS1=SS1-(LOG(RND))*((K+A+B)-INT(K+A+B))
220 FOR L1=1 TO INT(V+A-k):SS2=SS2-LOG(RND):NEXT
230 SS2=SS2-(LOG(RND))*((V+A-k)-INT(V+A-k))
240 X(I,J)=SS1/(SS2+SS1):NEXT:NEXT
250 FOR I=1 TO M:FOR J=1 TO N;PRINT USING "
#.#####";X(I,J);:NEXT:PRINT :NEXT

```

إن نتيجة تنفيذ البرنامج أعلاه عند قيم $\rho = 2/3$ و $m = 4$ و $\beta = 2$ و $\alpha = 1$ و $n = 10$ مثلًا ، هي المشاهدات المولدة التالية ، التي تخضع لنطوزع بيتا ذي الاربع متغيرات :

0.4601	0.7486	0.8640	0.2959	0.4480	0.5358	0.6322	0.8494	0.2820
0.7333	0.3218	0.8000	0.5599	0.6305	0.3129	0.6097	0.8640	0.7744
0.3868	0.5131	0.4664	0.9860	0.7510	0.3344	0.8437	0.5272	0.8524
0.9045	0.6982	0.4870	0.6001	0.8169	0.2160	0.4756	0.7563	0.5507
				0.7026	0.7920	0.3807	0.4606	

(24)

المصادر

- 1) Johnson, N. and Kotz, S. (1978) "Distributions in statistics : continuous multivariate Distributions" , wiley series , USA.
- 2) Michael, J. and Schucany, W. (2002) "The mixture approach for simulating bivariate Distributions with specific correlations" , The American Statisticians , 56 , pp.48-54 .
- 3) Minhajuddin, A. & Harris, I. and Schucany, W. (2004) "Simulating multivariate Distributions with specific correlations" , The American Statisticians , 58 , pp.86-98 .
- 4) Mood, A. & Graybill, F. and Boes, D. (1989) "Introduction to the theory of statistics" , third edition , McGraw Hill company , USA .
- 5) Patil , G. & Boswell , M. & Joshi , S. and Ratnaparkhi , M. (1984) "Dictionary of statistical distributions" , McGraw Hill company , USA .